

# The Natural Gradient by Analogy to Signal Whitening, and Recipes and Tricks for its Use

Jascha Sohl-Dickstein  
Redwood Center for Theoretical Neuroscience  
University of California at Berkeley

March 3, 2013

The natural gradient, as introduced by [Amari, 1987], allows for more efficient gradient descent by removing dependencies and biases inherent in a function’s parameterization. Several papers present the topic thoroughly and precisely [Amari, 1987, Amari, 1998, Amari and Nagaoka, 2000, Theis, 2005, Amari, 2010]. It remains a very difficult idea to get your head around however. The intent of this note is to provide simple intuition for the natural gradient and its uses. We review how an ill conditioned parameter space can undermine learning, introduce the natural gradient by analogy to the more widely understood concept of signal whitening, and present tricks and specific prescriptions for applying the natural gradient to learning problems. To our knowledge, this is the first time a connection has been made between signal whitening and the natural gradient.

## 1 Natural gradient

### 1.1 A simple example

We begin with a simple probabilistic model which has clearly been very poorly parametrized. For this we use a two dimensional gaussian distribution, with means written in terms of the parameters  $\theta \in \mathcal{R}^2$ ,

$$q(\mathbf{x}; \theta) = \frac{1}{2\pi} \exp \left[ -\frac{1}{2} \left( x_1 - \left[ 3\theta_1 + \frac{1}{3}\theta_2 \right] \right)^2 - \frac{1}{2} \left( x_2 - \left[ \frac{1}{3}\theta_1 \right] \right)^2 \right]. \quad (1)$$

As an objective function  $J(\theta)$  we use the negative log likelihood of  $q(\mathbf{x}; \theta)$  under an observed data distribution  $p(\mathbf{x})$

$$J(\theta) = -\langle \log q(\mathbf{x}; \theta) \rangle_{p(\mathbf{x})}. \quad (2)$$

Using steepest gradient descent to minimize the negative log likelihood involves taking steps like

$$\Delta\theta \propto -\nabla_{\theta}J(\theta) \quad (3)$$

$$\begin{bmatrix} \Delta\theta_1 \\ \Delta\theta_2 \end{bmatrix} \propto \begin{bmatrix} \langle 3(x_1 - [3\theta_1 + \frac{1}{3}\theta_2]) + \frac{1}{3}(x_2 - [\frac{1}{3}\theta_1]) \rangle_{p(\mathbf{x})} \\ \langle \frac{1}{3}(x_1 - [3\theta_1 + \frac{1}{3}\theta_2]) \rangle_{p(\mathbf{x})} \end{bmatrix}. \quad (4)$$

As can be seen in Figure 1a the steepest gradient update steps can move the parameters in a direction nearly perpendicular to the desired direction.  $q(\mathbf{x};\theta)$  is much more sensitive to changes in  $\theta_1$  than  $\theta_2$ , so the step size in  $\theta_1$  should be much smaller, but is instead much larger. In addition,  $\theta_1$  and  $\theta_2$  are not independent of each other. They move the distribution in nearly the same direction, making movement in the perpendicular direction particularly difficult. Getting the parameters here to fully converge via steepest descent is a slow proposition, as shown in Figure 1b.

The pathological learning gradient above is illustrative of a more general problem. A model's learning gradient is effected by the parameterization of the model as well as the objective function being minimized. The effects of the parameterization can dominate learning. The natural gradient is a technique to remove the effects of model parameterization from learning updates.

## 1.2 A metric on the parameter space

As a first step towards compensating for differences in relative scaling, and cross-parameter dependencies, the shape of the parameter space  $\theta$  is first described by assigning it a measure of distance, or a metric. This metric is expressed via a symmetric matrix  $\mathbf{G}(\theta)$ , which defines the length  $|d\theta|$  of an infinitesimal step  $d\theta$  in the parameters,

$$|d\theta|^2 = \sum_i \sum_j G_{ij}(\theta) d\theta_i d\theta_j = d\theta^T \mathbf{G}(\theta) d\theta. \quad (5)$$

$\mathbf{G}(\theta)$  is chosen so that the length  $|d\theta|$  provides a reasonable measure for the expected magnitude of the difference of  $J(\theta + d\theta)$  from  $J(\theta)$ . That is,  $\mathbf{G}(\theta)$  is chosen such that  $|d\theta|$  is representative of the expected magnitude of the change in the objective function resulting from a step  $d\theta$ . There is no uniquely correct choice for  $\mathbf{G}(\theta)$ .

If the objective function  $J(\theta)$  is the log likelihood of a probability distribution  $q(\mathbf{x};\theta)$ , then a measure of the information distance between  $q(\mathbf{x};\theta + d\theta)$  and  $q(\mathbf{x};\theta)$  usually works well, and the Fisher information matrix (Equation 30) is frequently used as a metric. Plugging in the example from Section 1.1, the resulting Fisher information matrix is  $\mathbf{G} = \begin{bmatrix} 3^2 + \frac{1}{3^2} & 1 \\ 1 & \frac{1}{3^2} \end{bmatrix}$ .

## 1.3 Connection to covariance

$\mathbf{G}(\theta)$  is an analogue of the inverse covariance matrix  $\Sigma^{-1}$ . Just as a signal can be whitened given  $\Sigma^{-1}$  — removing all first order dependencies and scaling

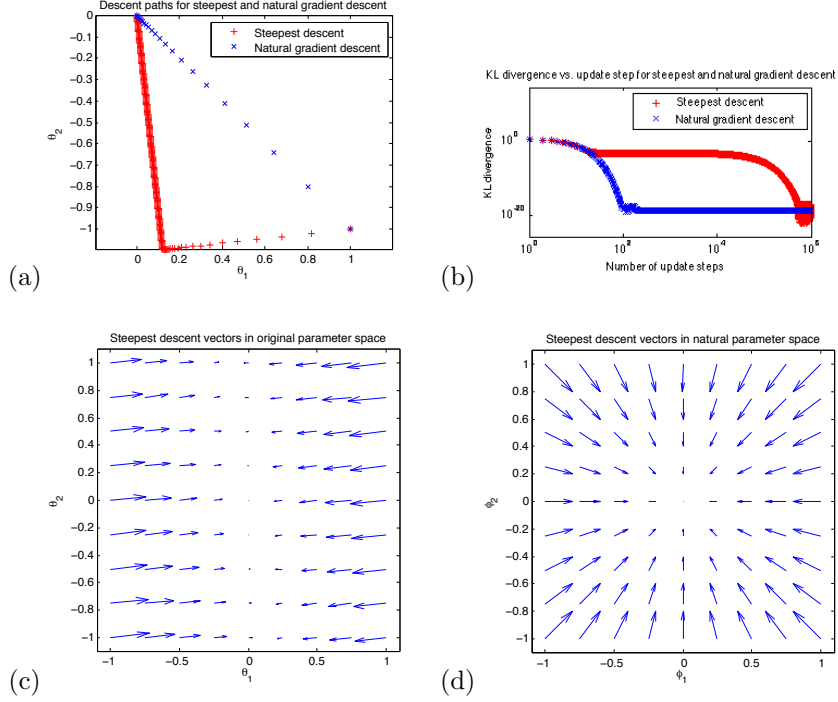


Figure 1: (a) The parameter descent paths taken by steepest gradient descent (red) and natural gradient descent (blue) for the example given in Section 1.1. The parameters are initialized at  $\theta_{init} = [1, -1]^T$ , and are fit to data generated with  $\theta_{true} = [0, 0]^T$ . The Fisher information matrix (Equation 30) is used to calculate the natural gradient. Notice that steepest descent takes a more circuitous and far slower path. (b) The KL divergence between the data distribution and the fit model as a function of number of gradient descent steps. Descent using the natural gradient converges more quickly. (c) The arrows give the gradient of the log likelihood objective (Equation 2), for a grid of parameter settings. This is the descent direction provided by Equation 4. (d) The gradient of the same log likelihood objective (Equation 2), but in terms of the whitened, natural, parameter space  $\phi$  as described in Section 1.4. Note that steepest descent in the whitened space converges directly to the true parameter values  $\phi_{true} = \mathbf{G}^{\frac{1}{2}}\theta_{true} = [0, 0]^T$ .

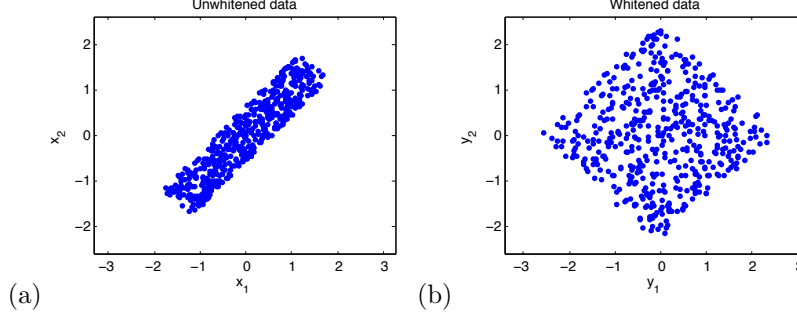


Figure 2: Example of signal whitening. (a) Samples  $\mathbf{x}$  from an unwhitened distribution in 2 variables. (b) The same samples after whitening, in new variables  $\mathbf{y} = \mathbf{W}\mathbf{x} = \mathbf{\Sigma}^{-\frac{1}{2}}\mathbf{x}$ .

the variance in each dimension to unit length — the parameterization of  $J(\theta)$  can also be “whitened,” removing the dependencies and differences in scaling between dimensions captured by  $\mathbf{G}(\theta)$ . See Figure 2 for an example of signal whitening.

As a quick review, the covariance matrix  $\mathbf{\Sigma}$  of a signal  $\mathbf{x}$  is defined as

$$\mathbf{\Sigma} = \langle \mathbf{x}\mathbf{x}^T \rangle. \quad (6)$$

The inverse covariance matrix is frequently used as a metric on the signal  $\mathbf{x}$ . This is called the Mahalanobis distance. It has the same form as the definition of  $|d\theta|^2$  in Equation 5,

$$|d\mathbf{x}|_{\text{Mahalanobis}}^2 = d\mathbf{x}^T \mathbf{\Sigma}^{-1} d\mathbf{x}. \quad (7)$$

In order to whiten a signal  $\mathbf{x}$ , a whitening matrix  $\mathbf{W}$  is found such that the covariance matrix for a new signal  $\mathbf{y} = \mathbf{W}\mathbf{x}$  is the identity matrix  $\mathbf{I}$ . The signal  $\mathbf{y}$  is then a whitened version of  $\mathbf{x}$ ,

$$\mathbf{I} = \langle \mathbf{y}\mathbf{y}^T \rangle = \mathbf{W} \langle \mathbf{x}\mathbf{x}^T \rangle \mathbf{W}^T = \mathbf{W}\mathbf{\Sigma}\mathbf{W}^T. \quad (8)$$

Remembering that  $\mathbf{\Sigma}^{-1}$  is symmetric, one solution<sup>1</sup> to this system of linear equations is

$$\mathbf{W} = \mathbf{\Sigma}^{-\frac{1}{2}} \quad (9)$$

$$\mathbf{y} = \mathbf{\Sigma}^{-\frac{1}{2}} \mathbf{x}. \quad (10)$$

If the covariance matrix for  $\mathbf{y}$  is the identity, then the metric for the Mahalanobis distance in the new variables  $\mathbf{y}$  is also the identity ( $|d\mathbf{y}|_{\text{Mahalanobis}}^2 = \mathbf{y}^T \mathbf{y}$ ).

---

<sup>1</sup> Choosing  $\mathbf{W} = \mathbf{\Sigma}^{-\frac{1}{2}}$  leads to symmetric, or zero-phase, whitening. In some fields it is referred to as a decorrelation stretch. It is equivalent to rotating a signal to the PCA basis, rescaling each axis to have unit norm, and then performing the inverse rotation, returning the signal to its original orientation. All unitary transformations of  $\mathbf{\Sigma}^{-\frac{1}{2}}$  also whiten the signal.

Whitening is a common preprocessing step in signal processing. It prevents incidental differences in scaling between dimensions from effecting later processing stages.

#### 1.4 “Whitening” the parameter space

If  $\mathbf{G}$  is not a function of  $\theta$ , then a similar procedure can be followed to produce a “whitened” parameterization  $\phi$ . We wish to find new parameters  $\phi = \mathbf{W}\theta$  such that the metric  $\mathbf{G}$  on  $\phi$  is the identity  $\mathbf{I}$ , as the Mahalanobis metric  $\mathbf{\Sigma}^{-1}$  is the identity for a whitened signal. This will mean that a small step  $d\phi$  in any direction will tend to have the same magnitude effect on the objective  $J(\phi)$ .

$$\phi = \mathbf{W}\theta \quad (11)$$

$$|d\phi|^2 = |d\theta|^2 \quad (12)$$

$$d\phi^T \mathbf{I} d\phi = d\theta^T \mathbf{G} d\theta \quad (13)$$

$$d\phi^T d\phi = d\theta^T \mathbf{G} d\theta \quad (14)$$

$$d\phi = \mathbf{W} d\theta \quad (15)$$

$$d\theta^T \mathbf{W}^T \mathbf{W} d\theta = d\theta^T \mathbf{G} d\theta \quad (16)$$

Noting that  $\mathbf{G}$  is symmetric, we find that one solution to this system of linear equations is

$$\mathbf{W} = \mathbf{G}^{\frac{1}{2}} \quad (17)$$

$$\phi = \mathbf{G}^{\frac{1}{2}} \theta \quad (18)$$

Steepest gradient descent steps in terms of  $\phi$  descend the objective function in a more direct fashion than steepest gradient descent steps in terms of  $\theta$ , as is illustrated in Figure 1c and 1d. In  $\phi$ , the steepest gradient is the natural gradient.

$\mathbf{G}$  is almost always a function of  $\theta$ , and for most problems there is no parameterization  $\phi$  which will be “white” everywhere. So long as  $\mathbf{G}(\theta)$  changes slowly though, it can be treated as constant for a single learning step. This suggests the following as an algorithm for learning in a natural parameter space.

1. Express  $J(\cdot)$  in terms of natural parameters  $\phi = \mathbf{G}^{\frac{1}{2}}(\theta_t) \theta$ .
2. Calculate an update step  $\Delta\phi \propto \nabla_{\phi} J(\phi_t)$ , where  $\phi_t = \mathbf{G}^{\frac{1}{2}}(\theta_t) \theta_t$ .
3. Calculate the  $\theta_{t+1} = \mathbf{G}^{-\frac{1}{2}}(\theta_t) (\phi_t + \Delta\phi)$  associated with the update to  $\phi$ .
4. Repeat.<sup>2</sup>

The resulting update steps more directly and rapidly descend the objective function than steepest descent steps.

---

<sup>2</sup>Practically,  $\mathbf{G}(\theta)$  can usually be treated as constant for many learning steps. This allows the natural gradient to be combined in a plug and play fashion with other gradient descent algorithms, like L-BFGS, by performing gradient descent on  $J(\phi)$  rather than  $J(\theta)$ .

## 1.5 The natural gradient in $\theta$

The parameter updates in Section 1.4 can be performed entirely in the original parameter space  $\theta$ . The natural gradient  $\tilde{\nabla}_\theta J(\theta)$  is the direction in  $\theta$  which is equivalent to steepest gradient descent in  $\phi$  of  $J(\phi)$ . In order to find  $\tilde{\nabla}_\theta J(\theta)$ , we first write  $\Delta\phi$  in terms of  $\theta$ , then we write the natural gradient update step in  $\theta$ ,  $\tilde{\Delta}\theta$ , in terms of  $\Delta\phi$ ,

$$\Delta\phi \propto \nabla_\phi J(\phi) \quad (19)$$

$$= \left( \frac{\partial\theta}{\partial\phi^T} \right)^T \nabla_\theta J(\theta) \quad (20)$$

$$= \mathbf{G}^{-\frac{1}{2}} \nabla_\theta J(\theta) \quad (21)$$

(where  $\frac{\partial\theta}{\partial\phi^T}$  is the Jacobian matrix),

$$\tilde{\Delta}\theta \propto \frac{\partial\theta}{\partial\phi^T} \Delta\phi \quad (22)$$

$$= \mathbf{G}^{-\frac{1}{2}} \Delta\phi \quad (23)$$

$$\propto \mathbf{G}^{-1} \nabla_\theta J(\theta). \quad (24)$$

Since the natural gradient update step is proportional to the natural gradient,  $\tilde{\Delta}\theta \propto \tilde{\nabla}_\theta J(\theta)$ , the natural gradient can be written as

$$\tilde{\nabla}_\theta J(\theta) = \mathbf{G}^{-1}(\theta) \nabla_\theta J(\theta) \quad (25)$$

Figure 1a illustrates this gradient applied to the example objective function from Section 1.1. If gradient descent is performed by infinitesimal steps in the direction indicated by  $\tilde{\nabla}_\theta J(\theta)$ , then the parameterization of the problem will have no effect on the path taken during learning (though choice of  $\mathbf{G}(\theta)$  will have an effect).

## 2 Recipes and tricks

In this section we present a reference with key formulas for using the natural gradient, as well as approaches useful for applying the natural gradient in specific cases.

### 2.1 Natural gradient

The natural gradient is

$$\tilde{\nabla}_\theta J(\theta) = \mathbf{G}^{-1}(\theta) \nabla_\theta J(\theta) \quad (26)$$

where  $J(\theta)$  is an objective function to be minimized with parameters  $\theta$ , and  $\mathbf{G}(\theta)$  is a metric on the parameter space. Learning should be performed with

an update rule

$$\theta_{t+1} = \theta_t + \tilde{\Delta}\theta_t \quad (27)$$

$$\tilde{\Delta}\theta \propto -\tilde{\nabla}_{\theta} J(\theta) \quad (28)$$

with steps taken in the direction given by the natural gradient.

## 2.2 Metric $G(\theta)$

If the objective function  $J(\theta)$  is the negative log likelihood of a probabilistic model  $q(\mathbf{x}; \theta)$  under an observed data distribution  $p(\mathbf{x})$

$$J(\theta) = -\langle \log q(\mathbf{x}; \theta) \rangle_{p(\mathbf{x})} \quad (29)$$

then the Fisher information matrix

$$G_{ij}(\theta) = \left\langle \frac{\partial \log q(\mathbf{x}; \theta)}{\partial \theta_i} \frac{\partial \log q(\mathbf{x}; \theta)}{\partial \theta_j} \right\rangle_{q(\mathbf{x}; \theta)} \quad (30)$$

is a good metric to use.

If the objective function is *not* of the form given in Equation 29, and cannot be transformed into that form, then greater creativity is required. See Section 2.8 for some basic hints.

Remember, as will be discussed in Section 2.10, even if the metric you choose is approximate, it is still likely to accelerate convergence!

## 2.3 Fisher information over data distribution

The Fisher information matrix (Equation 30) requires averaging over the model distribution  $q(\mathbf{x}; \theta)$ . For some models this is very difficult to do. If that is the case, instead taking the average over the empirical data distribution  $p(\mathbf{x})$

$$G_{ij}(\theta) = \left\langle \frac{\partial \log q(\mathbf{x}; \theta)}{\partial \theta_i} \frac{\partial \log q(\mathbf{x}; \theta)}{\partial \theta_j} \right\rangle_{p(\mathbf{x})} \quad (31)$$

is frequently an effective alternative.

## 2.4 Energy approximation

Parameter estimation in a probabilistic model of the form

$$q(\mathbf{x}) = \frac{e^{-E(\mathbf{x}; \theta)}}{Z(\theta)} \quad (32)$$

is in general very difficult, since it requires working with the frequently intractable partition function integral  $Z(\theta) = \int e^{-E(\mathbf{x}; \theta)} d\mathbf{x}$ . There are a number of techniques which can provide approximate learning gradients (eg minimum

probability flow [Sohl-Dickstein et al., 2011b, Sohl-Dickstein et al., 2011a], contrastive divergence [Welling and Hinton, 2002, Hinton, 2002], score matching [Hyvärinen, 2005], mean field theory, and variational bayes [Tanaka, 1998, Kappen and Rodriguez, 1997, Jaakkola and Jordan, 1997, Haykin, 2008]). Turning those gradients into natural gradients is difficult though, as the Fisher information depends on the gradient of  $\log Z(\theta)$ . Practically, simply ignoring the  $\log Z(\theta)$  terms entirely and using a metric

$$G_{ij}(\theta) = \left\langle \frac{\partial E(\mathbf{x}; \theta)}{\partial \theta_i} \frac{\partial E(\mathbf{x}; \theta)}{\partial \theta_j} \right\rangle_{p(\mathbf{x})} \quad (33)$$

averaged over the data distribution works surprisingly well, and frequently greatly accelerates learning.

## 2.5 Diagonal approximation

$\mathbf{G}(\theta)$  is a square matrix of size  $N \times N$ , where  $N$  is the number of parameters in the vector  $\theta$ . For problems with large  $N$ ,  $\mathbf{G}^{-1}(\theta)$  can be impractically expensive to compute and apply. For almost all problems however, the natural gradient still improves convergence even when off-diagonal elements of  $\mathbf{G}(\theta)$  are neglected,

$$G_{ij}(\theta) = \delta_{ij} \left\langle \left( \frac{\partial \log q(\mathbf{x}; \theta)}{\partial \theta_i} \right)^2 \right\rangle_{q(\mathbf{x}; \theta)}, \quad (34)$$

making inversion and application cost  $O(N)$  to perform.

If the parameters can be divided up into several distinct classes (for instance the covariance matrix and means of a gaussian distribution), block diagonal forms may also be worth considering.

## 2.6 Regularization

Even if evaluating the full  $\mathbf{G}$  is easy for your problem, you may still find that  $\mathbf{G}^{-1}$  is ill conditioned<sup>3</sup>. Dealing with this — solving a set of linear equations subject to some regularization, rather than using an unstable matrix inverse — is an entire field of study in computer science. Here we give one simple plug and play technique, called stochastic robust approximation (Section 6.4.1 in [Boyd and Vandenberghe, 2004]), for regularizing the matrix inverse. If  $\mathbf{G}^{-1}$  is replaced with

$$\mathbf{G}_{reg}^{-1} = (\mathbf{G}^T \mathbf{G} + \epsilon \mathbf{I})^{-1} \mathbf{G}^T \quad (35)$$

---

<sup>3</sup>This is a general problem when taking matrix inverses. A matrix  $\mathbf{A}$  with random elements, or with noisy elements, will tend to have a few very very small eigenvalues. The eigenvalues of  $\mathbf{A}^{-1}$  are the inverses of the eigenvalues of  $\mathbf{A}$ .  $\mathbf{A}^{-1}$  will thus tend to have a few very very large eigenvalues, which will tend to make the elements of  $\mathbf{A}^{-1}$  very very large. Even worse, the eigenvalues and eigenvectors which most dominate  $\mathbf{A}^{-1}$  are those which were smallest, noisiest and least trustworthy in  $\mathbf{A}$ .



where  $\epsilon$  is some small constant (say 0.01), the matrix inverse will be much better behaved.

Alternatively, techniques such as ridge regression can be used to solve the linear equation

$$\mathbf{G}(\theta) \tilde{\nabla}_{\theta} J(\theta) = \nabla_{\theta} J(\theta) \quad (36)$$

for  $\tilde{\nabla}_{\theta} J(\theta)$ .

## 2.7 Combining the natural gradient with other techniques using the natural parameter space $\phi$

It can be useful to combine the natural gradient with other gradient descent techniques. Blindly replacing all gradients with natural gradients frequently causes problems (line search implementations, for instance, depend on the gradients they are passed being the true gradients of the function they are descending). For a fixed value of  $\mathbf{G}$  though there is a natural parameter space.

$$\phi = \mathbf{G}^{\frac{1}{2}}(\theta_{fixed}) \theta \quad (37)$$

in which the steepest gradient is the same as the natural gradient.

In order to easily combine the natural gradient with other gradient descent techniques, fix  $\theta_{fixed}$  to the initial value of  $\theta$  and perform gradient descent over  $\phi$  using any preferred algorithm. After a significant number of update steps convert back to  $\theta$ , update  $\theta_{fixed}$  to the new value of  $\theta$ , and continue gradient descent in the new  $\phi$  space.

## 2.8 Natural gradient of non-probabilistic models

The techniques presented here are not unique to probabilistic models. The natural gradient can be used in any context where a suitable metric can be written for the parameters. There are several approaches to writing an appropriate metric.

1. If the objective function is of a form

$$J(\theta) = \langle l(\mathbf{x}; \theta) \rangle_{p(\mathbf{x})} \quad (38)$$

where  $\langle \cdot \rangle_{p(\mathbf{x})}$  indicates averaging over some data distribution  $p(\mathbf{x})$ , then it is sensible to choose a metric based on

$$G_{ij}(\theta) = \left\langle \frac{\partial l(\mathbf{x}; \theta)}{\partial \theta_i} \frac{\partial l(\mathbf{x}; \theta)}{\partial \theta_j} \right\rangle_{p(\mathbf{x})} \quad (39)$$

2. Similarly, the penalty function can be treated as if it is the log likelihood of a probabilistic model, and the corresponding Fisher information matrix used.

For example, the task of minimizing an L2 penalty function  $\|\mathbf{y} - \mathbf{f}(\mathbf{x}; \theta)\|^2$  over observed pairs of data  $p(\mathbf{x}, \mathbf{y})$  can be made probabilistic. Imagine that the L2 penalty instead represents a conditional gaussian  $q(\mathbf{y}|\mathbf{x}; \theta) \propto \exp\left(-\|\mathbf{y} - \mathbf{f}(\mathbf{x}; \theta)\|^2\right)$  over  $\mathbf{y}$ , and use the observed marginal  $p(\mathbf{x})$  over  $\mathbf{x}$  to build a joint distribution  $q(\mathbf{x}, \mathbf{y}; \theta) = q(\mathbf{y}|\mathbf{x}; \theta) p(\mathbf{x})$ .<sup>4</sup> This generates the metric:

$$G_{ij}(\theta) = \left\langle \frac{\partial \log [q(\mathbf{y}|\mathbf{x}; \theta) p(\mathbf{x})]}{\partial \theta_i} \frac{\partial \log [q(\mathbf{y}|\mathbf{x}; \theta) p(\mathbf{x})]}{\partial \theta_j} \right\rangle_{q(\mathbf{y}|\mathbf{x}; \theta) p(\mathbf{x})} \quad (40)$$

$$= \left\langle \frac{\partial \log q(\mathbf{y}|\mathbf{x}; \theta)}{\partial \theta_i} \frac{\partial \log q(\mathbf{y}|\mathbf{x}; \theta)}{\partial \theta_j} \right\rangle_{q(\mathbf{y}|\mathbf{x}; \theta) p(\mathbf{x})} \quad (41)$$

3. Find a set of parameter transformations  $T(\theta)$  which you believe the distance measure  $|d\theta|$  should be invariant to, and then find a metric  $\mathbf{G}(\theta)$  such that this invariance holds. That is find  $\mathbf{G}(\theta)$  such that the following relationship holds for any invariant transformation  $T(\theta)$ ,

$$|(\theta + d\theta) - \theta|^2 = |T(\theta + d\theta) - T(\theta)|^2. \quad (42)$$

A special case of this approach involves functions parametrized by a matrix, as presented in the next section.

## 2.9 $\mathbf{W}^T \mathbf{W}$

As derived in [Amari, 1998], if a function depends on a (square, non-singular) matrix  $\mathbf{W}$ , it frequently aids learning a great deal to take

$$\tilde{\Delta} \mathbf{W}_{nat} \propto \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} \mathbf{W}^T \mathbf{W}. \quad (43)$$

The algebra leading to this rule is complex, but as discussed in the previous section it falls out of a demand that the distance measure  $|d\mathbf{W}|$  be invariant to a set of transformations applied to  $\mathbf{W}$ . In this case, those transformations are right multiplication by any (non-singular) matrix  $\mathbf{Y}$ .

$$d\theta^T \mathbf{G}(\theta) d\theta = (d\theta Y)^T \mathbf{G}(\theta Y) (d\theta Y) \quad (44)$$

## 2.10 What if my approximation of $\Delta\theta_{nat}$ is wrong?

For any positive definite  $\mathbf{H}$ , movement in a direction

$$\tilde{\Delta}\theta = \mathbf{H}\Delta\theta \quad (45)$$

---

<sup>4</sup>Amari [Amari, 1998] suggests using some uninformative model distribution  $q(\mathbf{x})$  over the inputs, such as a gaussian distribution, rather than taking  $p(\mathbf{x})$  from the data. Either approach will likely work well.

will descend the objective function. If the wrong  $\mathbf{H}$  is used, gradient descent is performed in a suboptimal way ... which is the problem when steepest gradient descent is used as well. Making an educated guess as to  $\mathbf{H}$  rarely makes things worse, and frequently helps a great deal.

## References

- [Amari and Nagaoka, 2000] Amari, S. and Nagaoka, H. (2000). *Methods of Information Geometry*, volume 191 of *Translations of Mathematical Monographs*. American Mathematical Society.
- [Amari, 1987] Amari, S.-I. (1987). *Differential Geometry in Statistical Inference*, volume 10 of *IMS Lecture Notes - Monograph Series*. Inst of Mathematical Statistic.
- [Amari, 1998] Amari, S.-I. (1998). Natural Gradient Works Efficiently in Learning. *Neural Computation*, 10(2):251–276.
- [Amari, 2010] Amari, S.-i. (2010). Information geometry in optimization, machine learning and statistical inference. *Frontiers of Electrical and Electronic Engineering in China*, 5(3):241–260.
- [Boyd and Vandenberghe, 2004] Boyd, S. P. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge Univ Press.
- [Haykin, 2008] Haykin, S. (2008). *Neural networks and learning machines; 3rd edition*. Prentice Hall.
- [Hinton, 2002] Hinton, G. E. (2002). Training products of experts by minimizing contrastive divergence. *Neural Computation*, 14(8):1771–1800.
- [Hyvärinen, 2005] Hyvärinen, A. (2005). Estimation of non-normalized statistical models using score matching. *Journal of Machine Learning Research*, 6:695–709.
- [Jaakkola and Jordan, 1997] Jaakkola, T. and Jordan, M. (1997). A variational approach to Bayesian logistic regression models and their extensions. *Proceedings of the sixth international workshop on artificial intelligence and statistics*.
- [Kappen and Rodriguez, 1997] Kappen, H. and Rodriguez, F. (1997). Mean field approach to learning in Boltzmann Machines. *Pattern Recognition Letters*.
- [Sohl-Dickstein et al., 2011a] Sohl-Dickstein, J., Battaglino, P., and DeWeese, M. (2011a). New Method for Parameter Estimation in Probabilistic Models: Minimum Probability Flow. *Physical Review Letters*, 107(22):11–14.
- [Sohl-Dickstein et al., 2011b] Sohl-Dickstein, J., Battaglino, P. B., and DeWeese, M. R. (2011b). Minimum Probability Flow Learning. *International Conference on Machine Learning*, 107(22):11–14.

- [Tanaka, 1998] Tanaka, T. (1998). Mean-field theory of Boltzmann machine learning. *Physical Review Letters E*.
- [Theis, 2005] Theis, F. (2005). Gradients on matrix manifolds and their chain rule. *Neural Information Processing-Letters and Reviews*.
- [Welling and Hinton, 2002] Welling, M. and Hinton, G. (2002). A new learning algorithm for mean field Boltzmann machines. *Lecture Notes in Computer Science*.